

WEEK 1:

Recurrent Neural Networks (RNNs)

Example:

$x:$ Word1 Word2 Word3 word4 — word9
 $x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad x^{<t>} \quad \dots \quad x^{<9>}$

$T_x = 9$ (length of input sequence)

$y:$ 1 1 0 1 1 0 0 0 0
 $y^{<1>} \quad y^{<2>} \quad \dots \quad \quad \quad y^{<9>}$

$T_y = 9$ (length of output sequence)

T_x and T_y need not be equal always.

$x^{(i)<t>}$

$T_x^{(i)}$

i^{th} example

length of sequence of i^{th} example

$y^{(i)<t>}$

$T_y^{(i)}$

Date:

Representing words :

X : Sentence of ~~different~~ words

Vocabulary

a	1	Usually size is
aaron	2	30K, 50K
:	:	
and	367	and even 100k
:	:	is not uncommon.
harry	4075	
:	:	
potter	6830	
:	:	
zulu	10000	

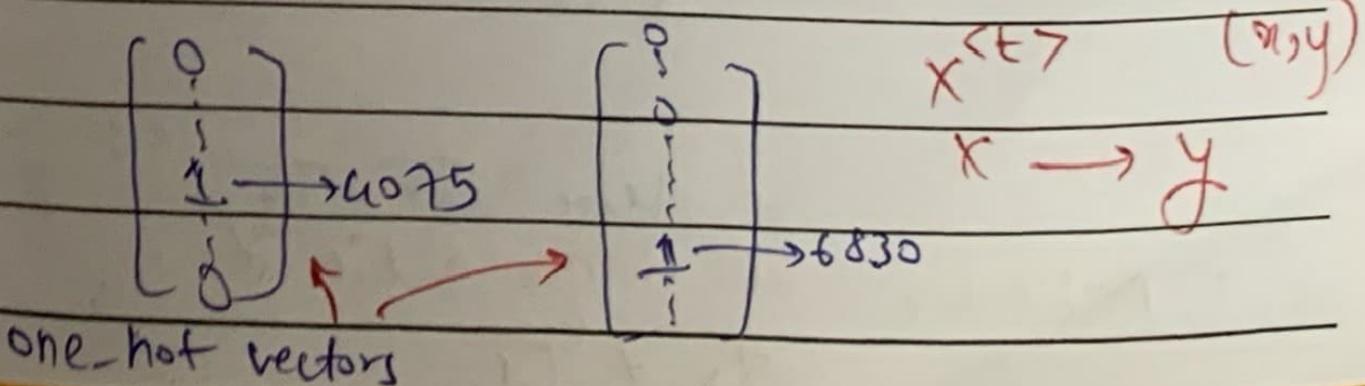
size: 10,000

Then using one-hot() all words can be assigned a vector where 1 is assigned

harry
 $x^{(1)}$

potter
 $x^{(2)}$

respective ^{to their} index

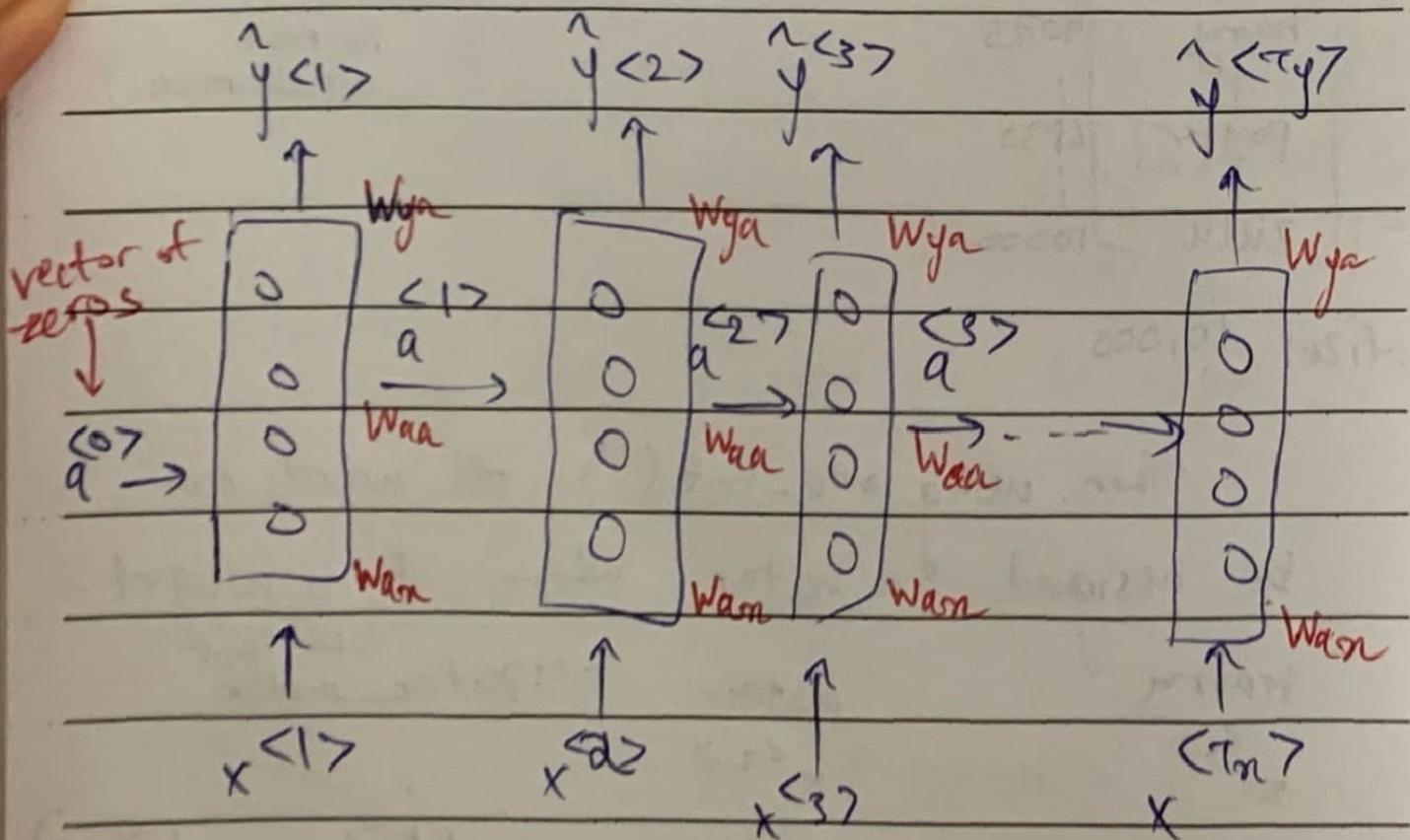


Date: _____

Why not a standard NN?

Problems:

- ① T_x and T_y will be different in every example, so doesn't work very well.
- ② Doesn't share features learned across different positions of text.



Here $T_x = T_y$
Architecture will change a bit if

$$T_x \neq T_y$$

Date:

→ The limitation of this NN is that to predict the output at the current time step, it uses info from the previous time steps but cannot use the later time steps.

This indicates that
it is used for computation of
" a " and will be multiplied by " y_n ".

$a^{(0)} = \vec{0}$ Most commonly used is $\tanh()$ and rarely ReLU	$a^{(1)} = g_1(W_a a^{(0)} + W_{\bar{a}n} \bar{x}^{(1)} + b_a)$
------------------------------------------------------------------------	-----------------------------------------------------------------

for binary classification
sigmoid otherwise softmax
but depends on requirement

$$a^{(t)} = g(W_a a^{(t-1)} + W_{\bar{a}n} \bar{x}^{(t)} + b_a)$$

$$\hat{y}^{(t)} = g(W_y a^{(t)} + b_y)$$

Date:

Simplified RNN notation:

$$a^{<t>} = g(w_a [a^{<t-1>}, x^{<t>}] + b_a)$$

$$y^{<t>} = g(w_y a^{<t>} + b_y)$$

If $w_{aa} : 100$

$w_{ax} : 10000$

then

$w_{aa} : 100, 100$

$w_{ax} : 100, 10000$

Stacked

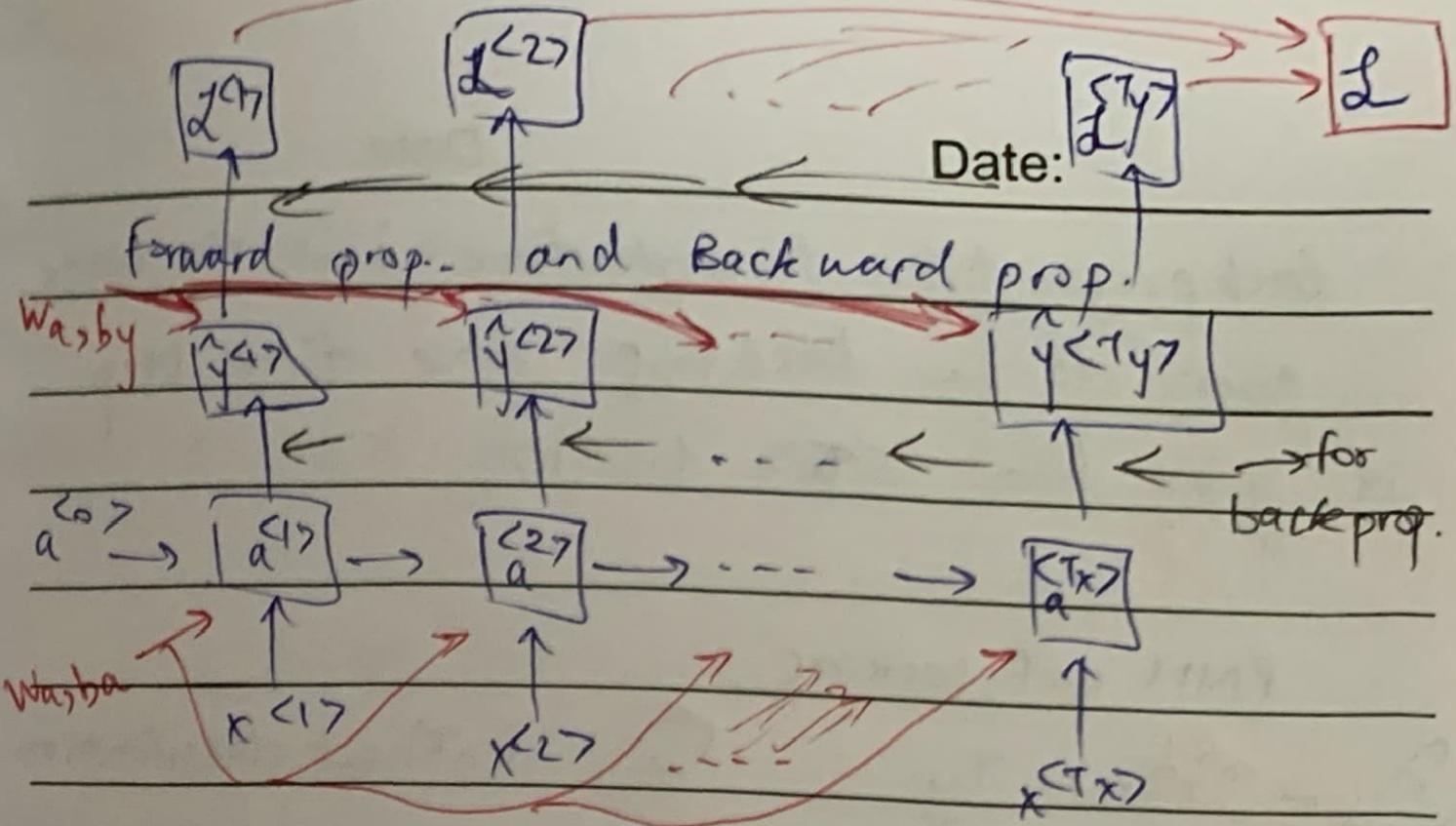
Stacked together as

$$\begin{bmatrix} w_{aa} & w_{ax} \end{bmatrix} = \begin{bmatrix} 100 & \\ & 10000 \end{bmatrix}$$

$\therefore w_a : (100, 10000)$

$$\begin{bmatrix} a^{<t-1>} & x^{<t>} \end{bmatrix} = \begin{bmatrix} a^{<t-1>} \\ \dots \\ x^{<t>} \end{bmatrix} \begin{matrix} \nearrow 100 \\ \times 10000 \\ \searrow 10100 \end{matrix}$$

$$\therefore \begin{bmatrix} w_{aa} & w_{ax} \end{bmatrix} \begin{bmatrix} a^{<t-1>} \\ \dots \\ x^{<t>} \end{bmatrix} = w_{aa} a^{<t-1>} + w_{ax} x^{<t>}$$



$$L(\hat{y}^{(t)}, y^{(t)}) = -\hat{y}^{(t)} \log \hat{y}^{(t)} - (1 - \hat{y}^{(t)}) \log (1 - \hat{y}^{(t)})$$

loss for single time step single word

$$L(\hat{y}, y) = \sum_{t=1}^{T_y} L^{(t)} (\hat{y}^{(t)}, y^{(t)})$$

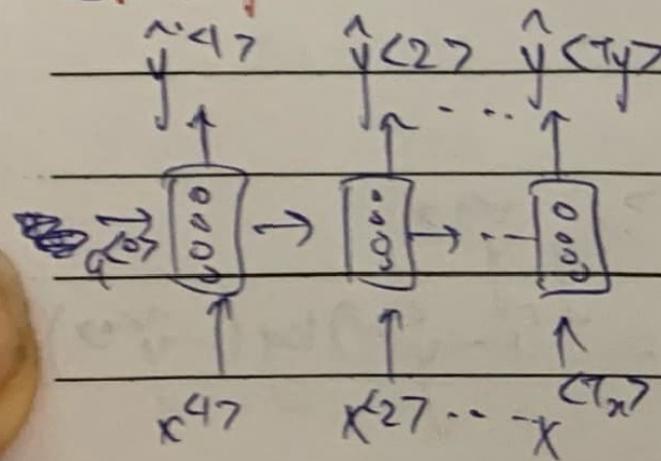
~~loss for entire sequence~~
~~loss for an~~
~~single word~~
~~loss for entire~~
~~sequence~~

Date: _____

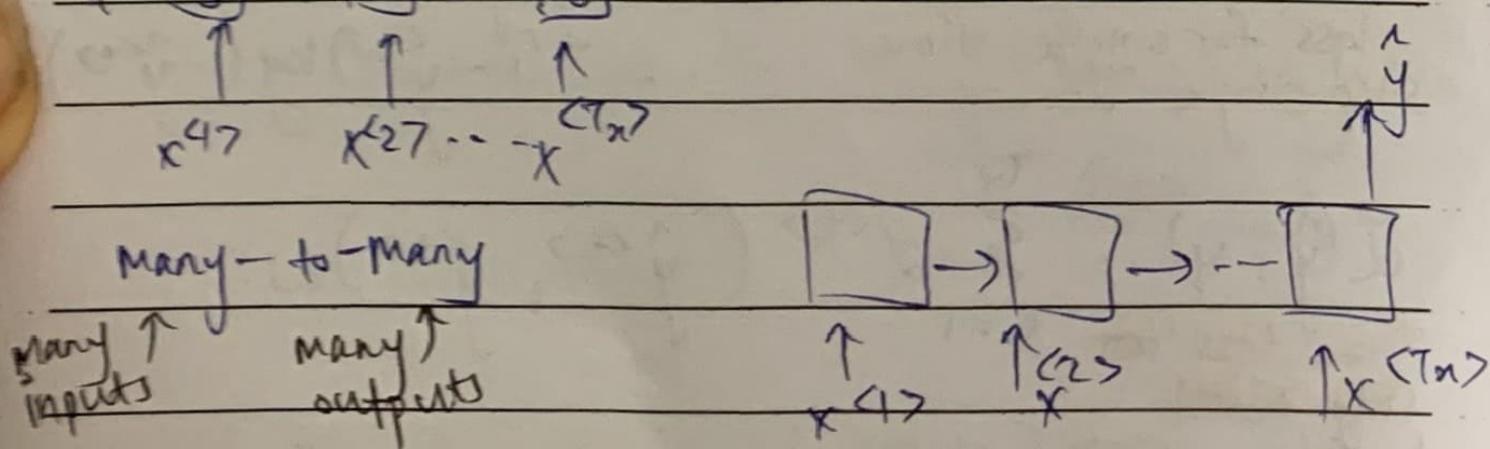
Backpropagation through time ; is the name given to the backprop since time step is going from $t+1$ to t

RNN architectures :

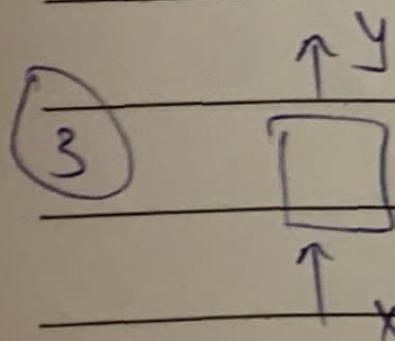
① Name entity recognition $x = T_y$ ② Sentiment classification



$x = \text{text sequence}$
 $y = 0/1 \text{ or } 1-5$



Many - to - One
Many inputs Many outputs
One number



One - to - One (similar almost same as standard NN)

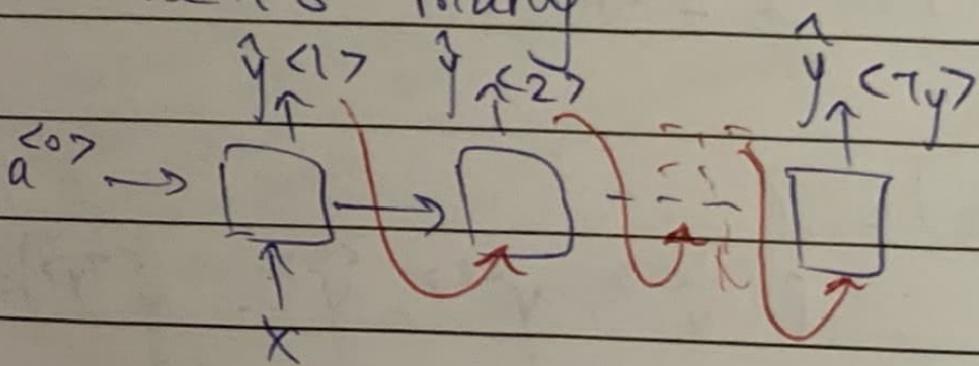
Date:

④ Music generation

$$x \rightarrow y^{<1>} y^{<2>} \dots y^{<T_y>}$$

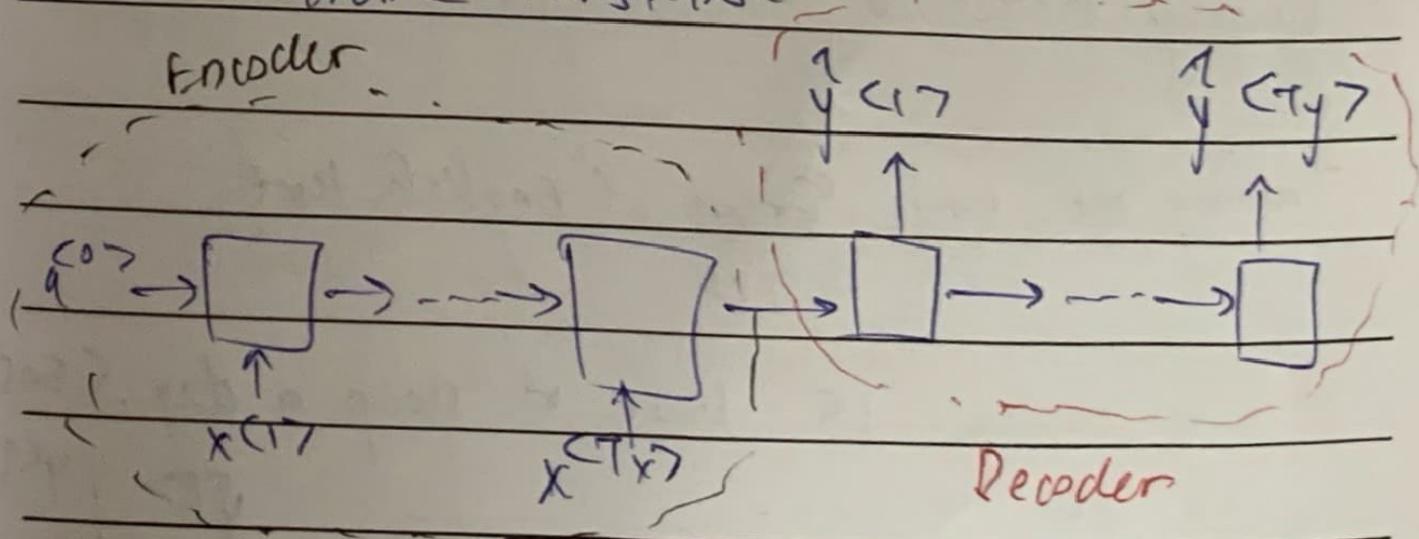
x can even be $x = \emptyset$

One-to-many



⑤ Machine translation

Encoder



Many-to-many $T_x \neq T_y$

Date: _____

What is language modelling?

→ Speech recognition:

The apple and air salad.

The apple and pear salad

$$P(1^{\text{st}}) = 3.2 \times 10^{-13}$$

$$P(2^{\text{nd}}) = 5.7 \times 10^{-10}$$

$$P(\text{sentence}) = ?$$

$$P(y^{<1>} y^{<2>} y^{<3>} \dots y^{<T_y>}) =$$

Language modelling with a RNN:

Training set: large ~~copy~~
corpus of English text.

Cats average 15 hours of sleep a day. (EOS)

$y^{<1>} y^{<2>} \dots y^{<8>} y^{<9>}$

↓
End of
sentence
token

$$P(y^{(1)}, y^{(2)}, y^{(3)}) = P(y^{(1)}) \cdot P(y^{(2)} | y^{(1)}) \cdot P(y^{(3)} | y^{(2)}, y^{(1)})$$

Date: [Redacted]

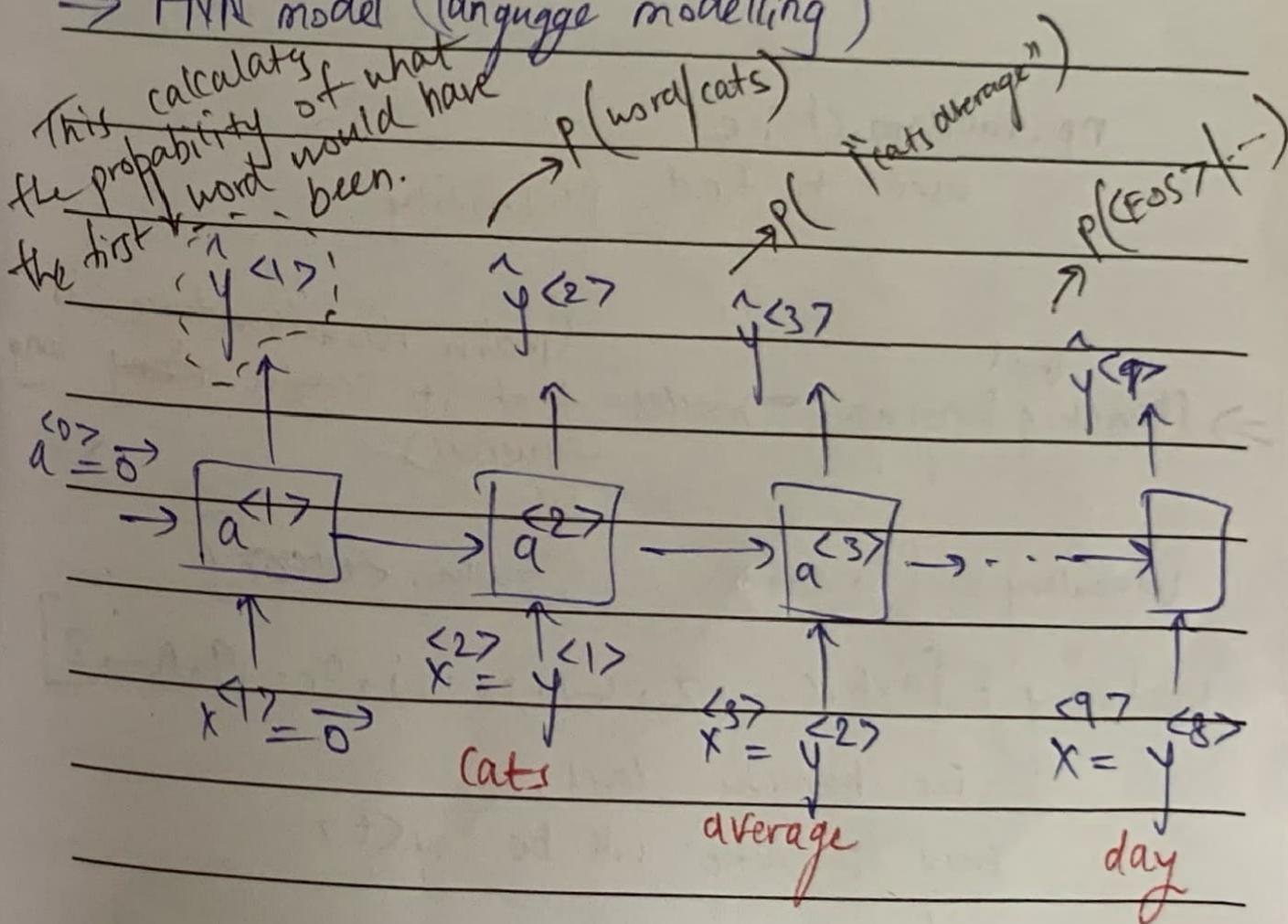
The Egyptian Main is a breed of cat. <EOS>

<UNK>

Unknown word

(not in the
dictionary-vocabulary
of 10k words used)

⇒ PNN model (language modelling)

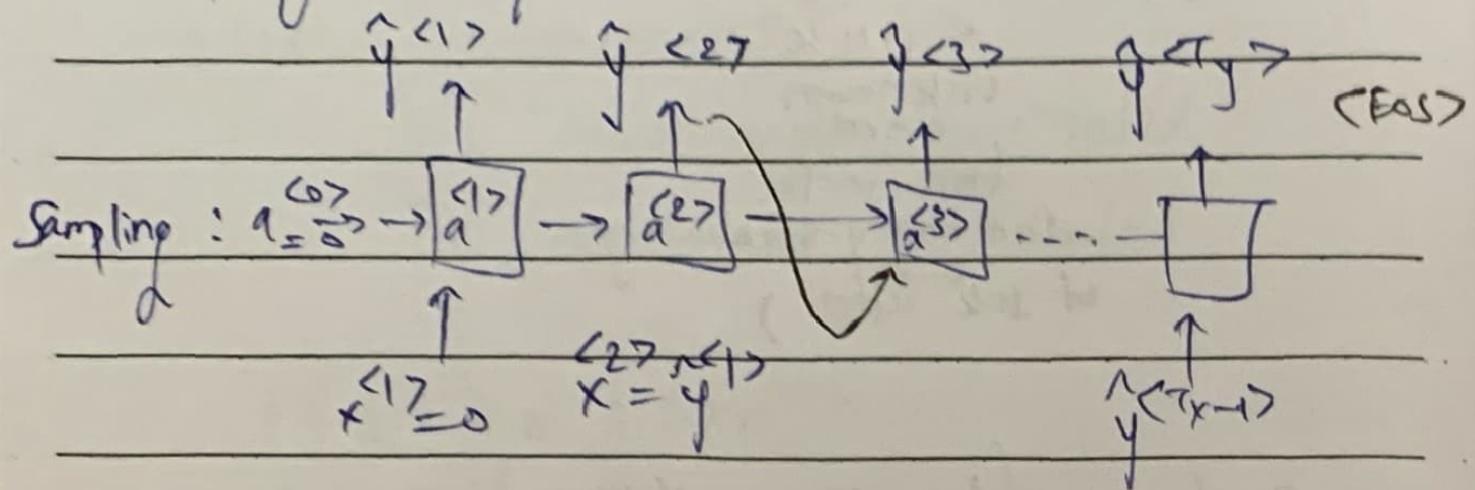


$$\mathcal{L}(y^{(t)}, \hat{y}^{(t)}) = - \sum_i y_i^{(t)} \log \hat{y}_i^{(t)}$$

$$\mathcal{L} = \sum_t \mathcal{L}^{(t)}(y^{(t)}, \hat{y}^{(t)})$$

Date:

Sampling a sequence from a trained RNN;



np.random.choice

used to find probability

⇒ character level language = model → (main disadvantage is that it creates ~~very~~ long sequences)

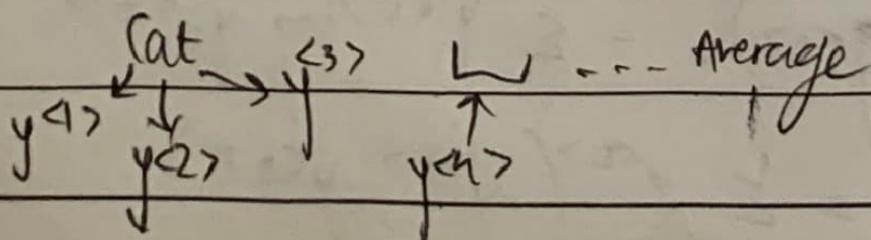
Vocabulary = [a, aaron, --, zulu, <UNK>]

Vocabulary = [a, b, c, --, z, l, o, i, j, 0, --, 9, A, Z]

for character level,

every character will be $y^{<t>}$

like :



Date:

Character level not good at long range dependencies. and more computationally expensive to train.

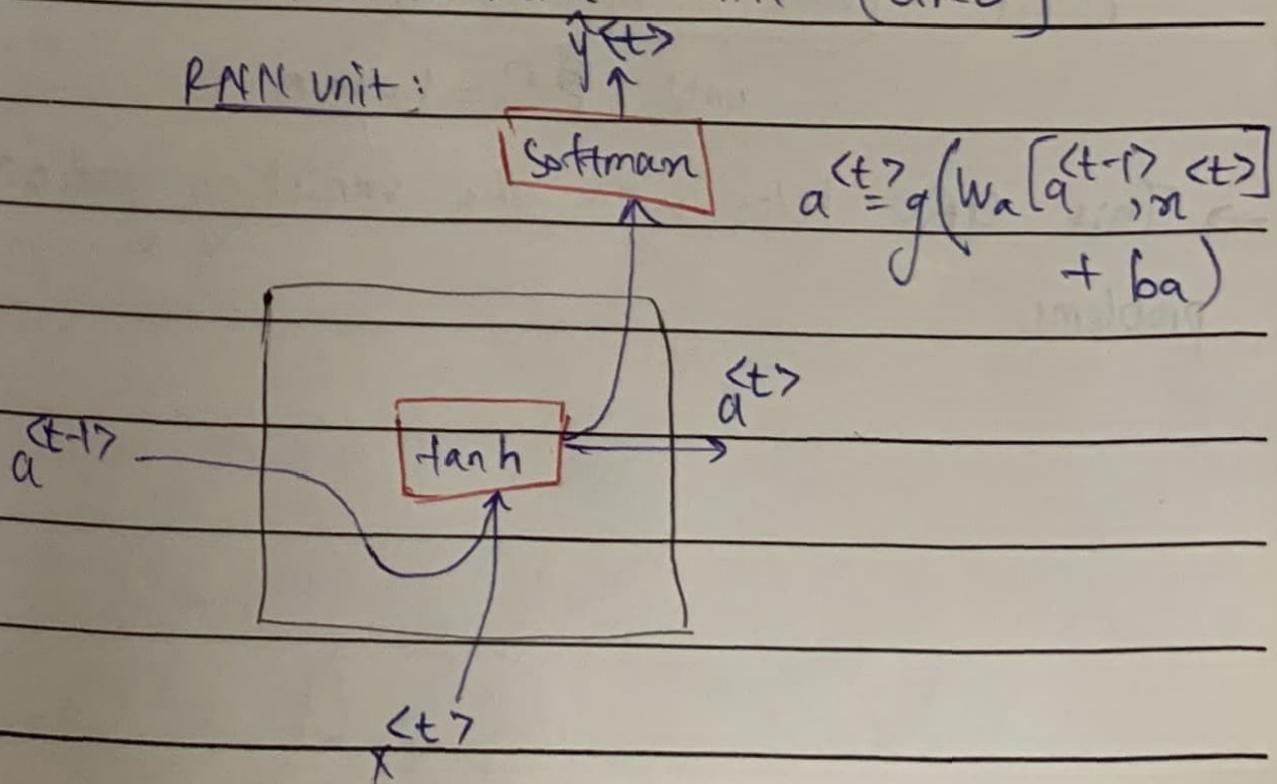
⇒ Vanishing gradients with RNNs:

The ^{basic} RNNs are not very good at capturing long term dependencies.

To manage Exploding gradients, use (gradient clipping)

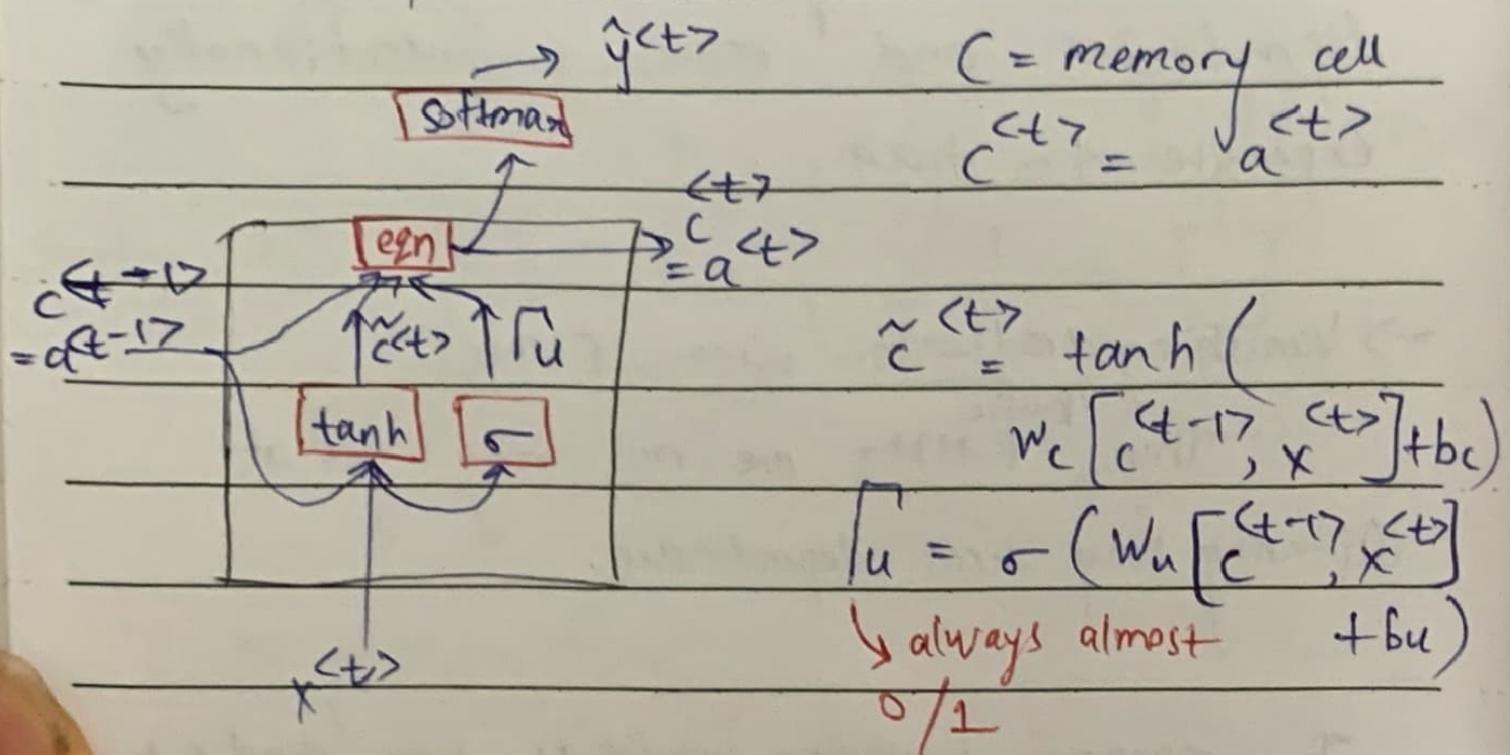
⇒ Gated Recurrent Unit [GRU]

RNN unit:



Γ_u = gamma u
 ↓
 "update"
 In egn,
 * is element-
 Date: wise mulp.

GRU (simplified) :



$$\text{eqn} \rightarrow c^{(t)} = \Gamma_u * \hat{c}^{(t)} + (1 - \Gamma_u)$$

\therefore if $\Gamma_u = 1$ $\Gamma_u = 0$ $\Gamma_u = 0 \dots \Gamma_u = 1$ $* c^{(t-1)}$

$c^{(t)} = 1$ (doesn't update)

until $\Gamma_u = 1$ again

\rightarrow This GRU helps with the vanishing gradient problems.

Date:

Full GRU:

$$\tilde{c}^{(t)} = \tanh(W_c \left[\underbrace{r^{(t-1)}}_{\text{in}} * \tilde{c}^{(t-1)}, x^{(t)} \right] + b_c)$$

$$r_u = \sigma(W_u \left[c^{(t-1)}, x^{(t)} \right] + b_u)$$

$$r_r = \sigma(W_r \left[c^{(t-1)}, x^{(t)} \right] + b_r)$$

$$c^{(t)} = r_u * \tilde{c}^{(t)} + (1 - r_u) * c^{(t-1)}$$

r = gate relevant tells how relevant the $c^{(t-1)}$ is to compute $\tilde{c}^{(t)}$

LSTM:

$$\tilde{c}^{(t)} = \tanh(W_c \left[\tilde{a}^{(t-1)}, x^{(t)} \right] + b_c)$$

"update" $r_u = \sigma(W_u \left[\tilde{a}^{(t-1)}, x^{(t)} \right] + b_u)$

"forget" $f_f = \sigma(W_f \left[\tilde{a}^{(t-1)}, x^{(t)} \right] + b_f)$

"output" $r_o = \sigma(W_o \left[\tilde{a}^{(t-1)}, x^{(t)} \right] + b_o)$

$$c^{(t)} = r_u * \tilde{c}^{(t)} + f_f * c^{(t-1)}$$

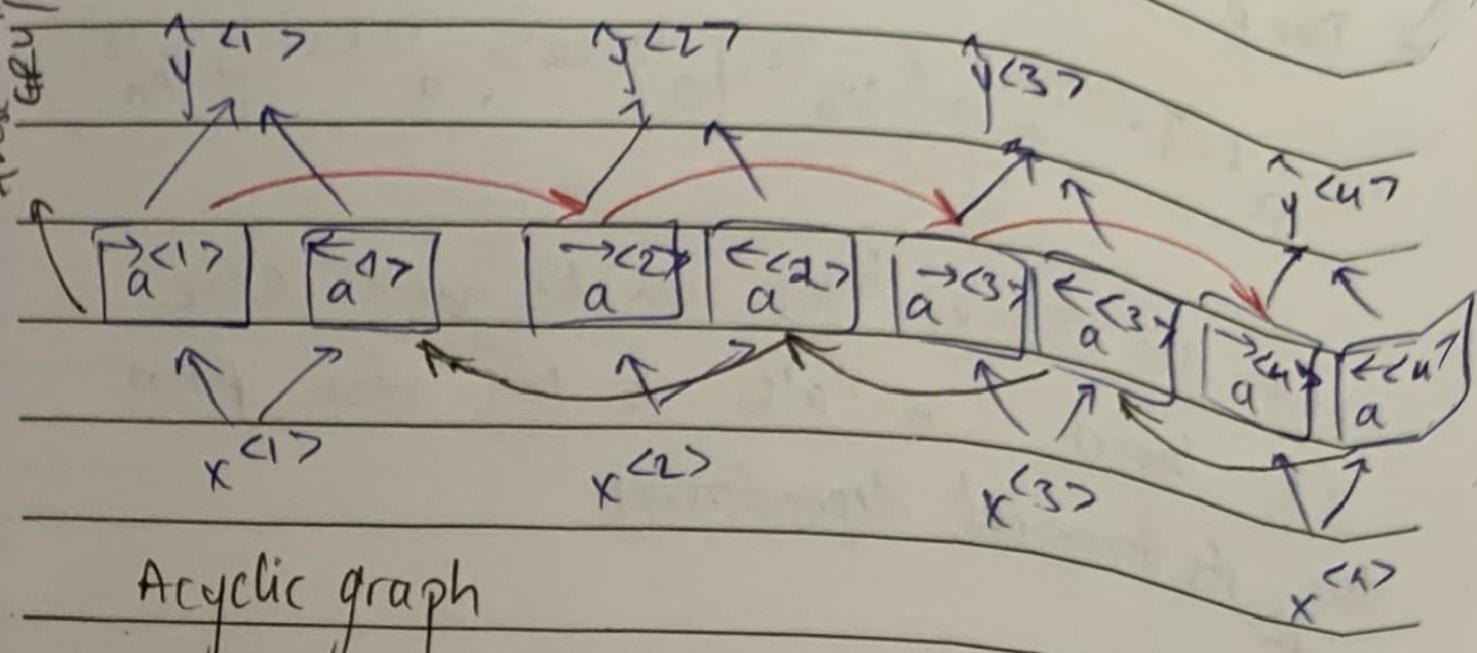
$$a^{(t)} = r_o * \tanh(c^{(t)})$$

\Rightarrow peep hole connection: variation of LSTM where
 c_{t-1} is also used to decide the value of
all three gate i_t , f_t , u_t Date:

LSTM in pictures:

Date:

⇒ Bidirectional RNN (BRNN):



$$\hat{y}^{(1)} = g(w_y [\vec{a}^{(t)}, \leftarrow a^{(t)}] + b_y)$$

This BRNN architecture works the whole sequence is to be present whole already at the start like NLP applications, whereas won't work for the real world speech applications.

These BRNNs can use tokens from all earlier and later part of the sequence.

Date:

Deep RNNs

$$a^{(2)} \leftarrow g\left(W_a^{(2)} \left[a^{(2)} \leftarrow, a^{(1)} \right] + b_a^{(2)}\right)$$

Not too many layers can be used since
3/4 layers is also a lot for RNN
due to its temporal dependencies.

